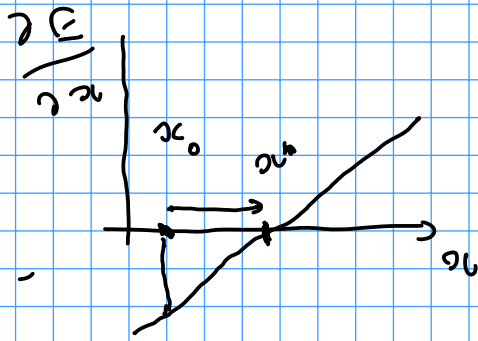


$$E(x) = \frac{1}{2}a(x - x^*)^2 + E(x^*)$$

$$\frac{\partial E(x)}{\partial x} = a(x - x^*)$$

$$\frac{\partial^2 E(x)}{\partial x^2} = a$$



$$\left(\frac{\frac{\partial E}{\partial x}(x_0)}{x_0 - x^*} \right) = a = \frac{\partial^2 E}{\partial x^2}(x_0)$$

$$x^* = x_0 - \left(\frac{\partial^2 E}{\partial x^2} \right)^{-1} \frac{\partial E}{\partial x}(x_0)$$

$$E(w) = \frac{1}{P} \sum_{p=1}^P \frac{1}{2} (y^p - w^T x^p)^2$$

$$\frac{\partial E}{\partial w} = -\frac{1}{P} \sum_{p=1}^P (y^p - w^T x^p) \cdot x^p$$

$$H = \frac{\partial^2 E}{\partial w \partial w^T} = \frac{1}{P} \sum_{p=1}^P x^p x^{pT}$$

Covariance matrix of samples

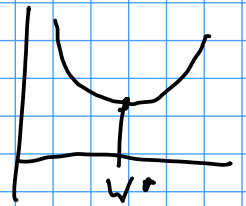
Hessian matrix

$$w_{t+1} = w_t - H(w_t)^{-1} \frac{\partial E}{\partial w}(w_t)$$

$$\left(\begin{array}{l} \text{Solve for } w_{t+1}: \\ H(w_{t+1}) (w_{t+1} - w_t) = \frac{\partial E}{\partial w}(w_t) \end{array} \right.$$

/

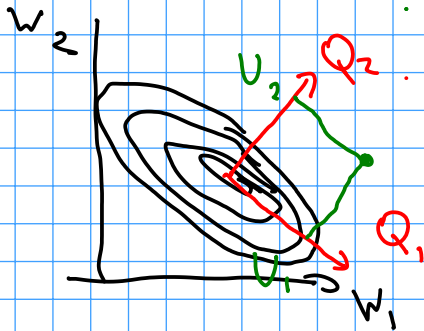
$$E(w) = E(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$



$$= + \frac{1}{2} (w - w^*)^T Q^T \Lambda Q (w - w^*)$$

$$\underline{v = Qw}$$

$$E(v) = E(v^*) + \frac{1}{2} (v - v^*)^T \Lambda (v - v^*)$$



$$v_{t+1} = v_t - \underbrace{\Lambda^{-1}}_{H^{-1}} \underbrace{\Lambda (v_t - v^*)}_{\text{gradient}} \rightarrow v_{t+1} = v^*$$

$$v_{t+1}(i) = v_t(i) - \underbrace{\lambda_i^{-1}}_{\text{step size}} \left[\frac{\partial E}{\partial v}(i) \right]$$

$$W_i \leftarrow W_i - \eta_i \left(\frac{\partial E}{\partial W} \right)_i$$

diagonal Newton

$$\eta_i = \frac{1}{\left[\frac{\partial^2 E}{(\partial W_i)^2} \right]^+ + \gamma}$$

Levenberg-Marquardt

$$(H + \gamma I)^{-1} \frac{\partial E}{\partial W}$$

positive semi-definite
estimate of Hessian

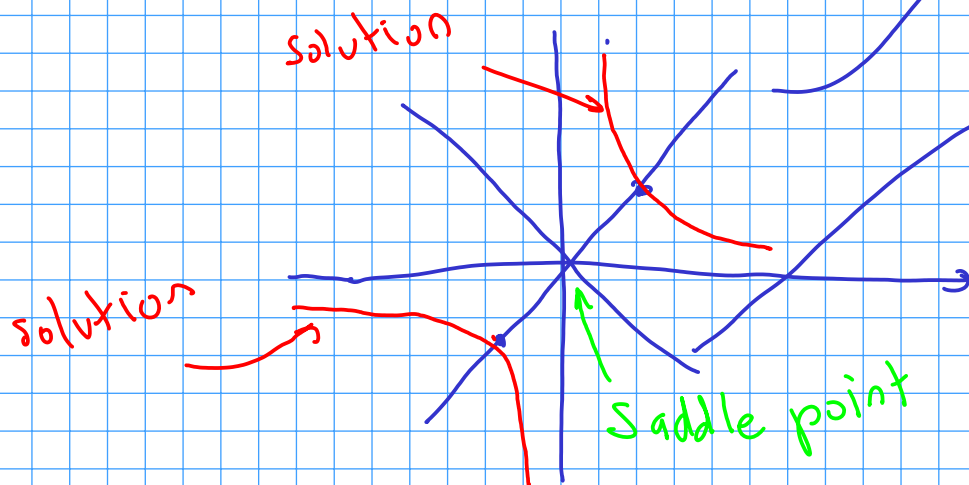
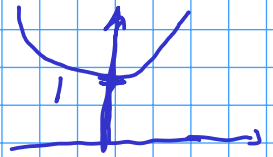
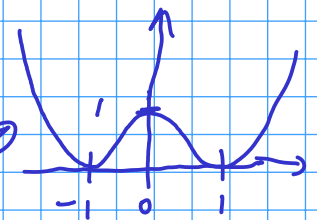
$$E(W) = (1 - W_1 W_2)^2$$

$$W = W_1 = W_2$$

$$W' = W_1 = -W_2$$

$$(1 - W^2)^2$$

$$(1 + W'^2)^2$$



$$M_{t+1} = M_t + \alpha \frac{\partial E}{\partial v}$$

Nesterov momentum

$$W_{t+1} = W_t - \eta M_t$$



$$H = \frac{1}{|p|} \sum_{p=1}^p x_p x_p^T$$

speed

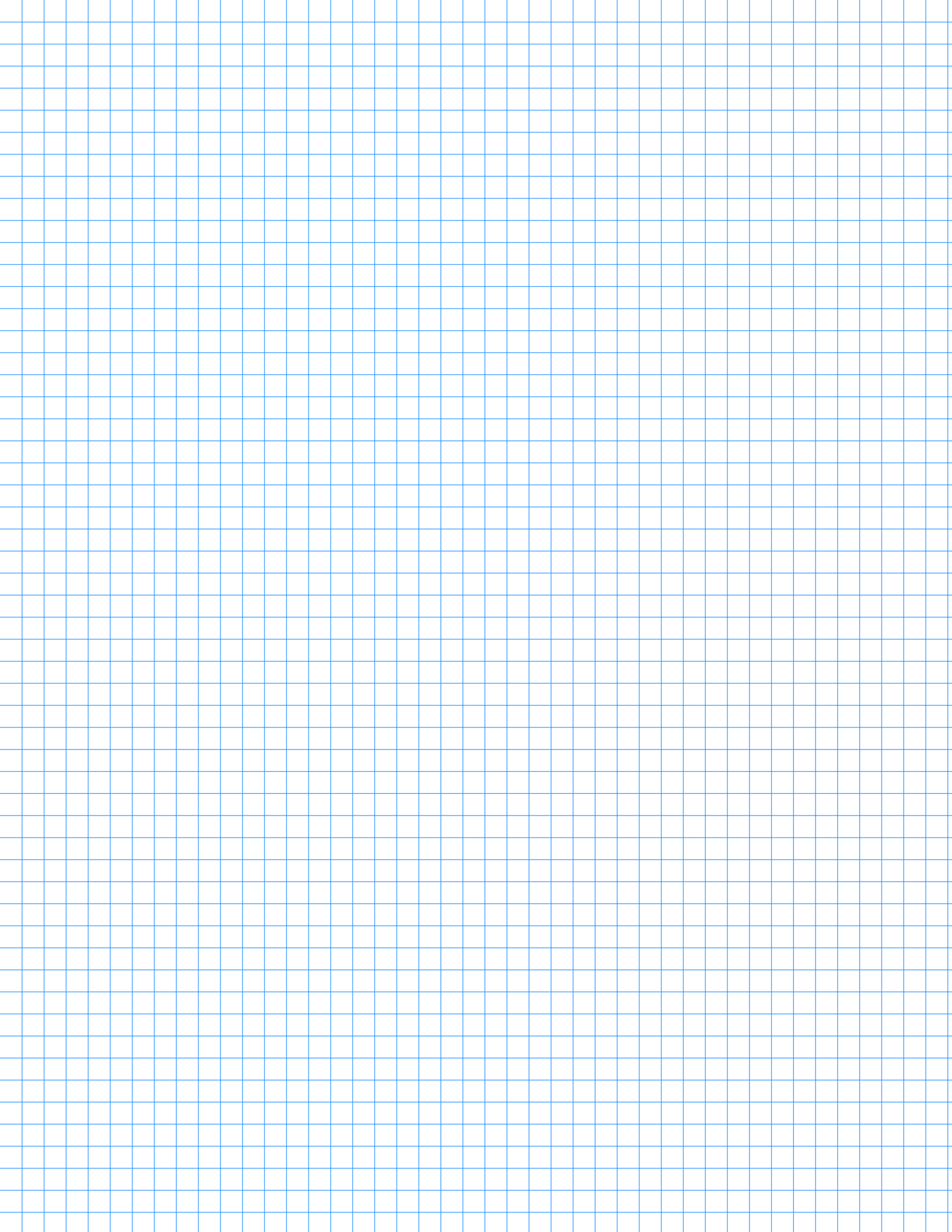
$$H = Q^T \Lambda Q$$

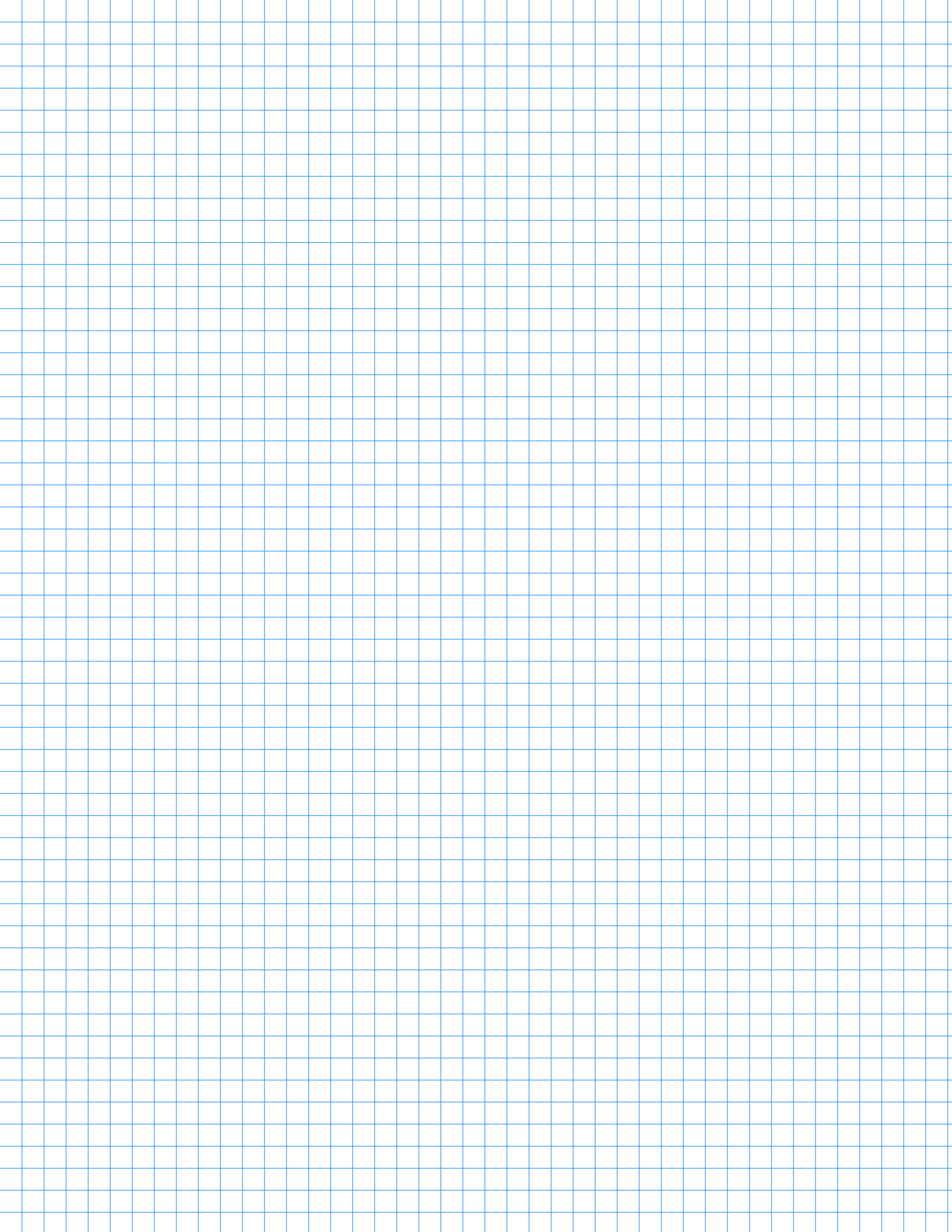
$$\Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix}$$



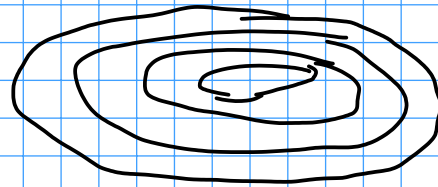
eigenvalues of H

eigenvectors of H

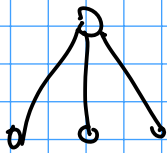
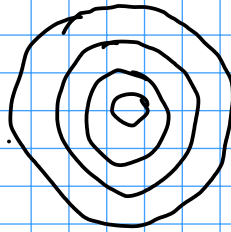




When H is diagonal
principal axes are aligned
with frame of reference



$$\begin{bmatrix} \lambda & & & \\ & \lambda & & \\ & & \lambda & \\ & & & \lambda \end{bmatrix}$$



$$H = \frac{1}{p} \sum_{i=1}^p x^p x^{p^T} = \lambda I$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix}$$

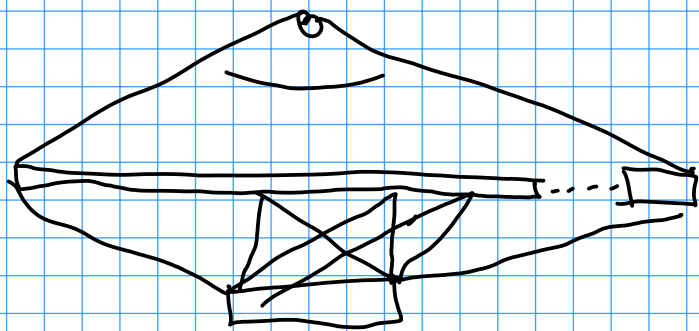
x_i :-
- zero mean
- variance 1
- decorrelated

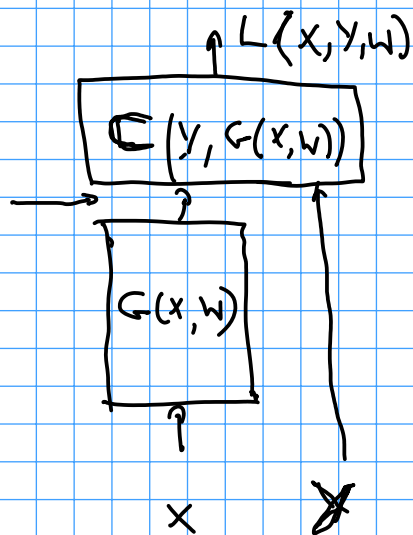
$$H_{ij} = \sum_{p=1}^p x_i^p x_j^p$$

$X^p \rightarrow Z^p$ 0 mean, variance 1, uncorrelated
└ whitening
Karchunen - Loeve transform

$$X \quad H = Q^T \Lambda Q \quad Z = \Lambda^{-1} Q (X - \bar{X})$$

$$H = \sum_p (x^p - \bar{x}^p)(x^p - \bar{x}^p)^T$$





Supervised learning

C differentiable

G differentiable

y observed

x observed

$L(x, y, w)$ computed



SGD Backprop

C non diff

G diff / G non diff

y obs

x obs

L computed



black-box optim

REINFORCE

RL: y is not observed

L is observed

C unknown

